# FPGA-accelerated machine learning inference as a solution for particle physics computing challenges

Jennifer Ngadiuba, Maurizio Pierini **(CERN)**
Javier Duarte, Burt Holzman, Ben Kreis, Kevin Pedro, *Mia Liu*, Nhan Tran, Aristeidis Tsaris **(Fermilab)**
Phil Harris, Dylan Rankin **(MIT)**
Zhenbin Wu **(UIC)**

**Dec 9, CPAD, 2018**

## Motivation
*Challenges of big science and computing*

## Our solution : proof of concept
*Particle physics computing with Brainwave*

## Physics cases
*NOvA & jet identification at collider experiments*

*Outlook& takeaways*

# Motivation: Challenges of big science and computing

## CMS as an example: Detectors becoming increasingly complex

- High-resolution detector
- Order of 100 Million channels



CMS DETECTOR
Total weight        : 14,000 tonnes
Overall diameter : 15.0 m
Overall length     : 28.7 m
Magnetic field     : 3.8 T

STEEL RETURN YOKE
12,500 tonnes

SILICON TRACKERS
Pixel (100x150 μm) ~16m² ~66M channels
Microstrips (80x180 μm) ~200m² ~9.6M channels

SUPERCONDUCTING SOLENOID
Niobium titanium coil carrying ~18,000A

MUON CHAMBERS
Barrel: 250 Drift Tube, 480 Resistive Plate Chambers
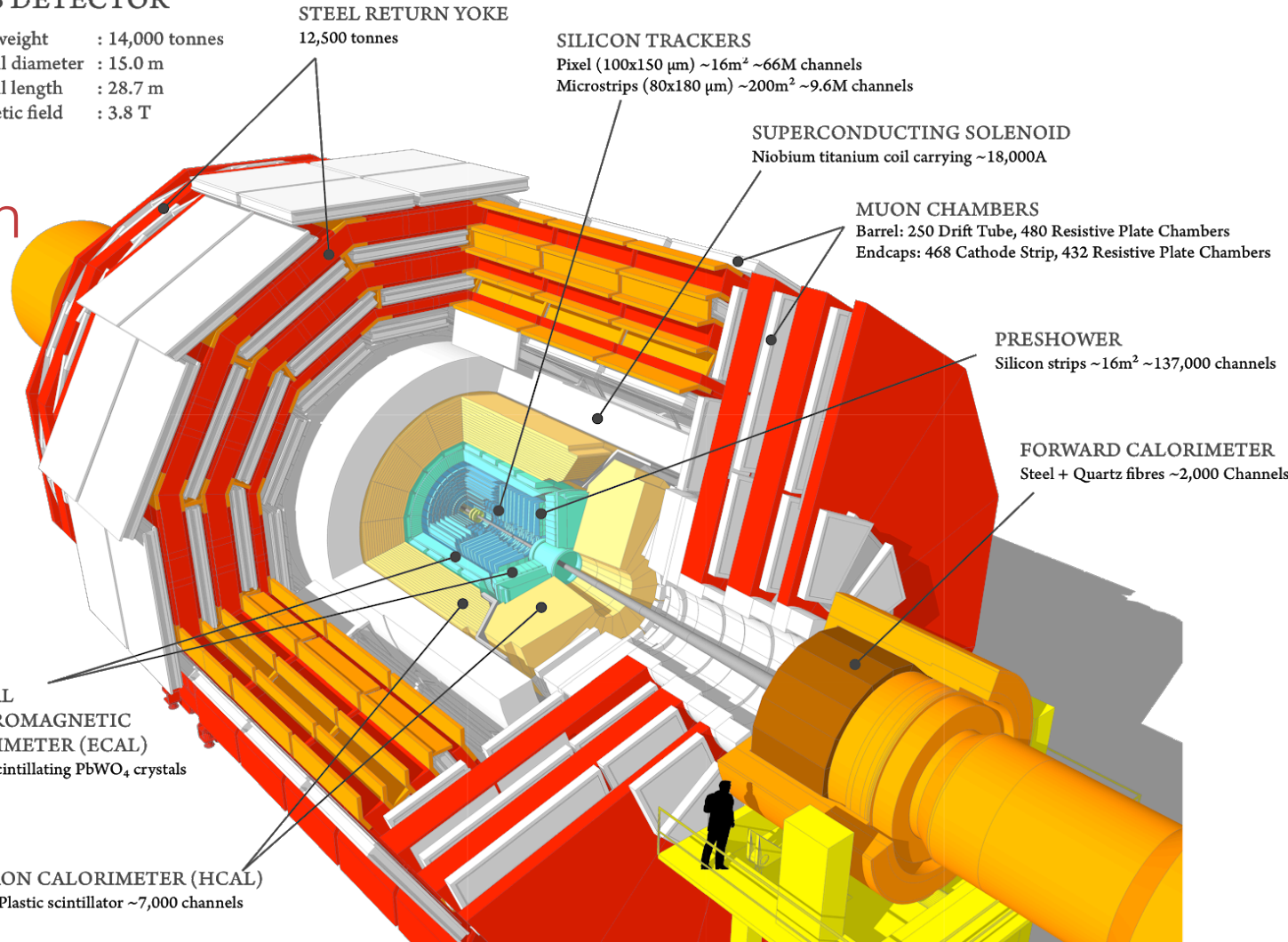Endcaps: 468 Cathode Strip, 432 Resistive Plate Chambers

PRESHOWER
Silicon strips ~16m² ~137,000 channels

FORWARD CALORIMETER
Steel + Quartz fibres ~2,000 Channels

CRYSTAL ELECTROMAGNETIC CALORIMETER (ECAL)
~76,000 scintillating PbWO₄ crystals

HADRON CALORIMETER (HCAL)
Brass + Plastic scintillator ~7,000 channels

# DETECTORS GET

**CMS upgrade to get**
**granularity, timing in**
**Exam... CMS High**

Total Silicon:

- 600 m$^2$

Total scintillator

- 500 m$^2$

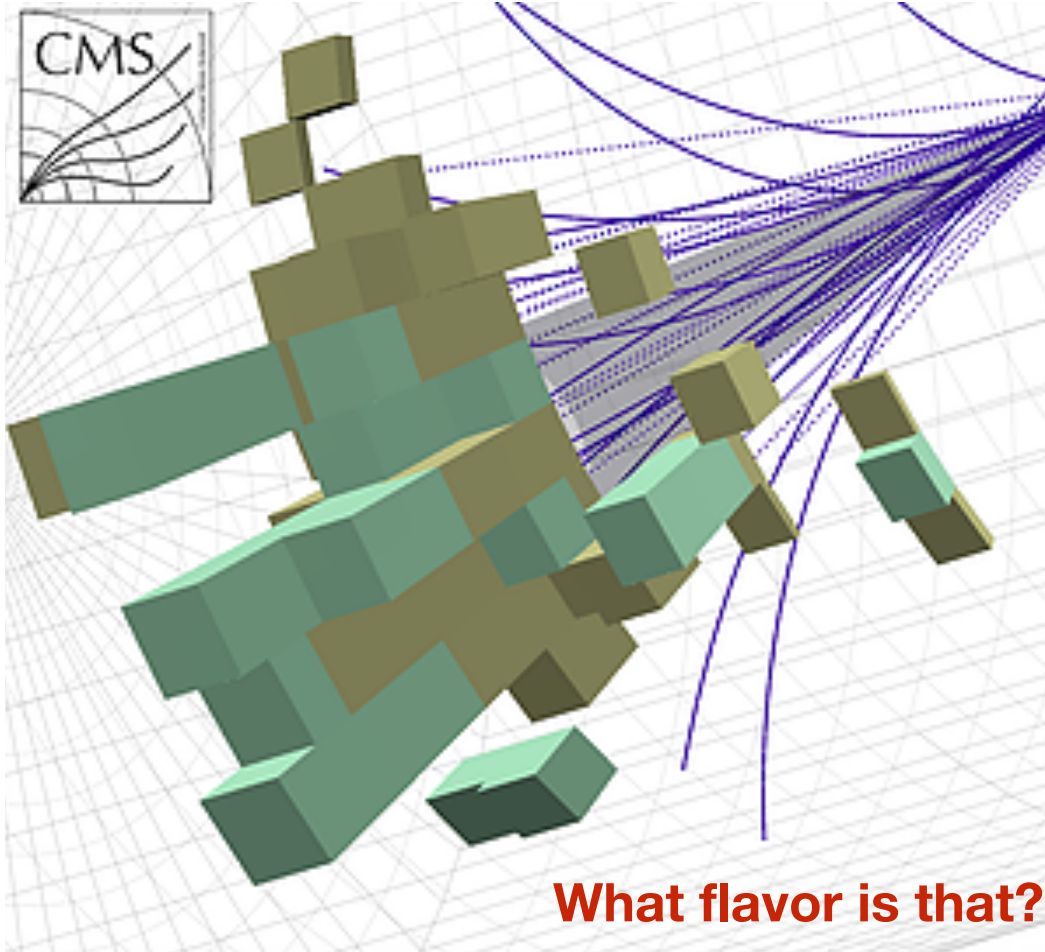| | CMS | ATLAS | CMS HGCal |
|---|---|---|---|
| Diameter (m) | 15 | 25 | |
| Length (m) | 28.7 | 46 | |
| B-Field (T) | 3.8 | 2/4 | |
| EM Cal channels | ~80,000 | ~110,000 | 4.3M |
| Had Cal channels | ~7,000 | ~10,000 | 1.8M |

P.Merkel

**CMS as an example: Need sophisticated algorithms to fully exploit the information taken by more complex detectors**
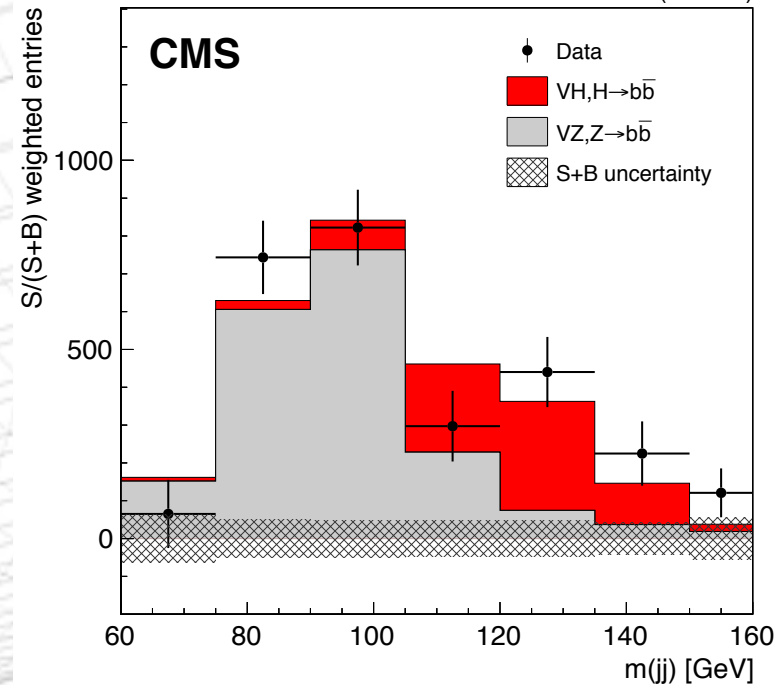


What flavor is that?

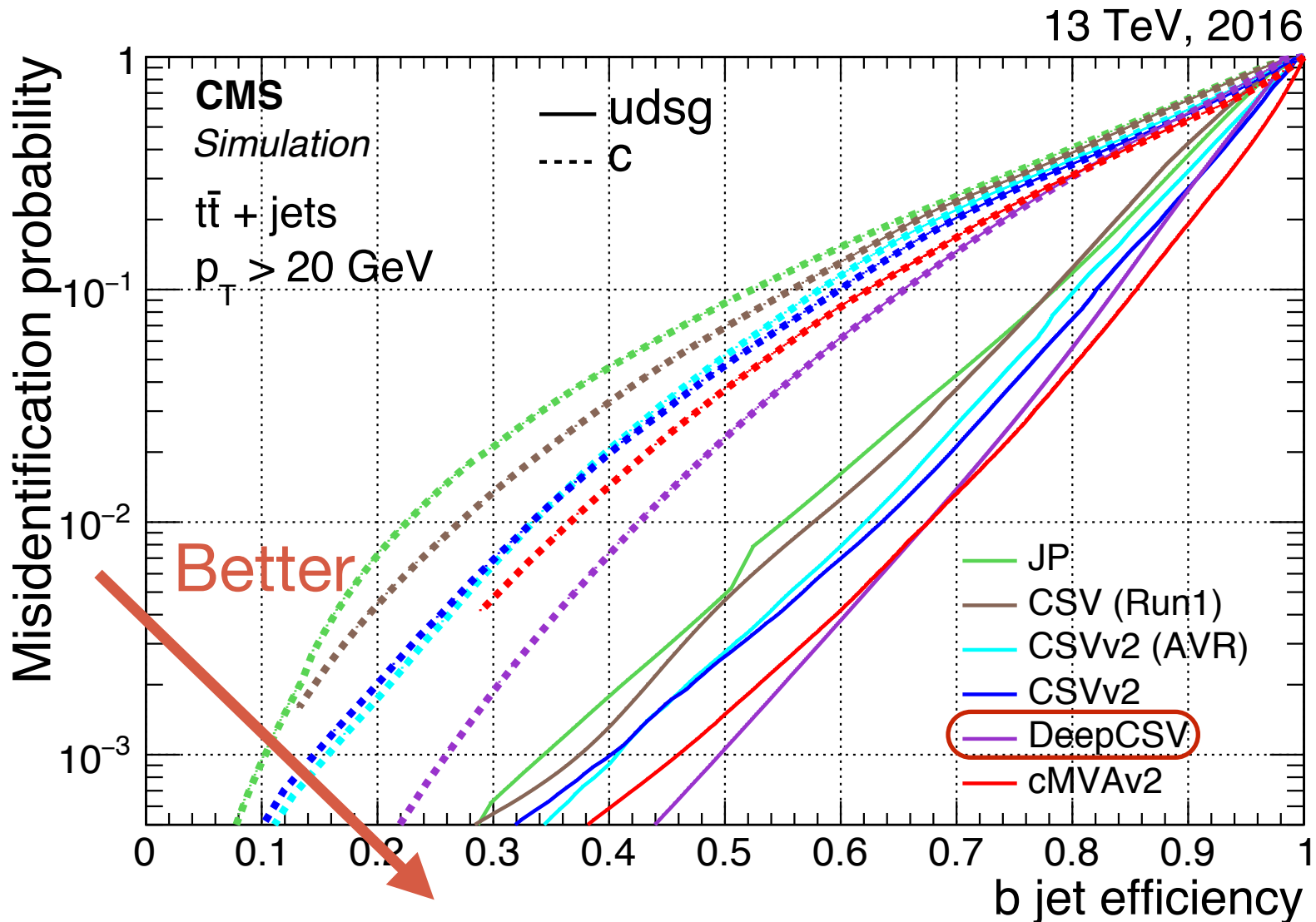**CMS as an example: plenty of physics cases**



Yukawa coupling: H->bb

What flavor is that?

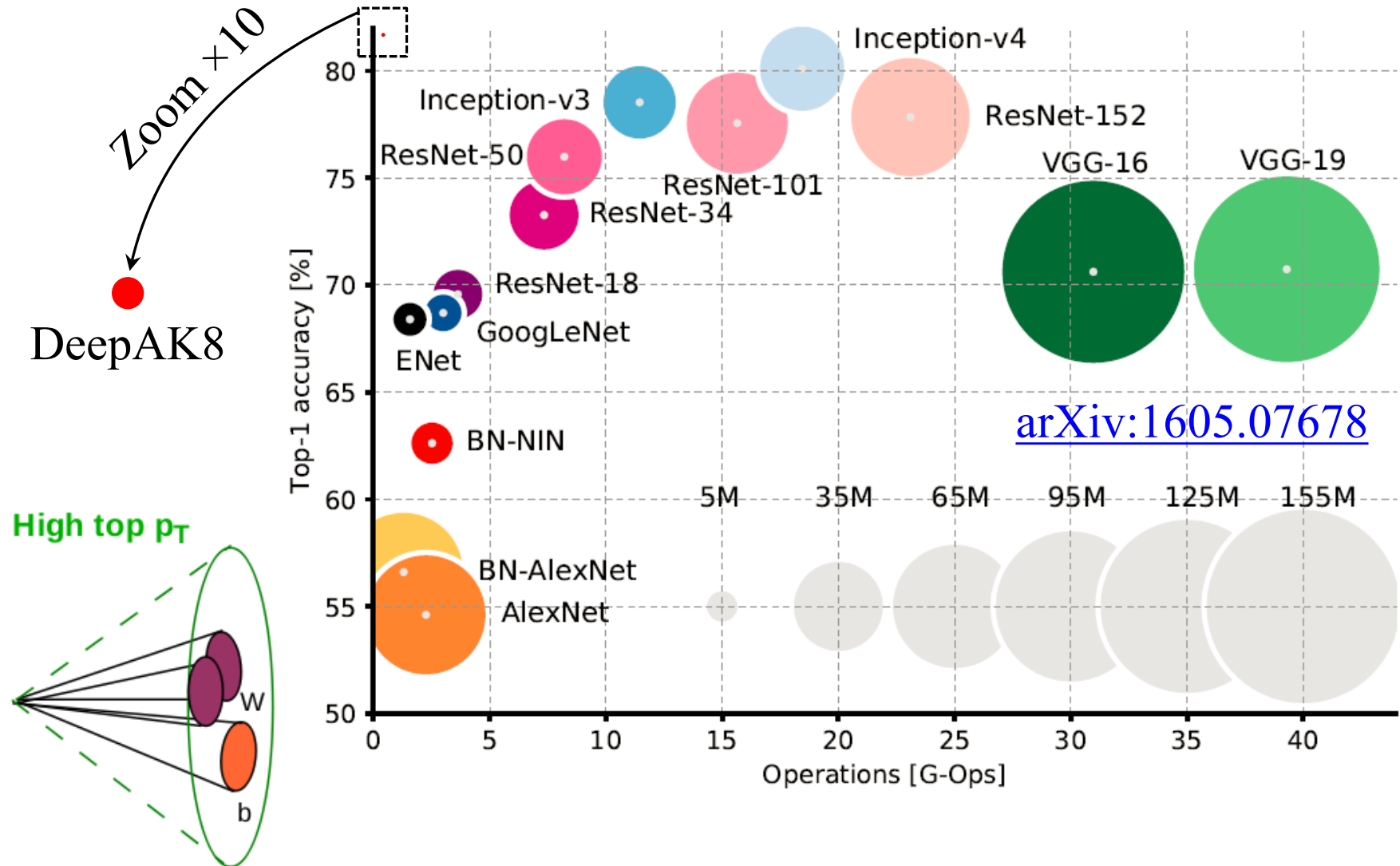# BOOM IN USING DEEP NEURAL NETWORKS

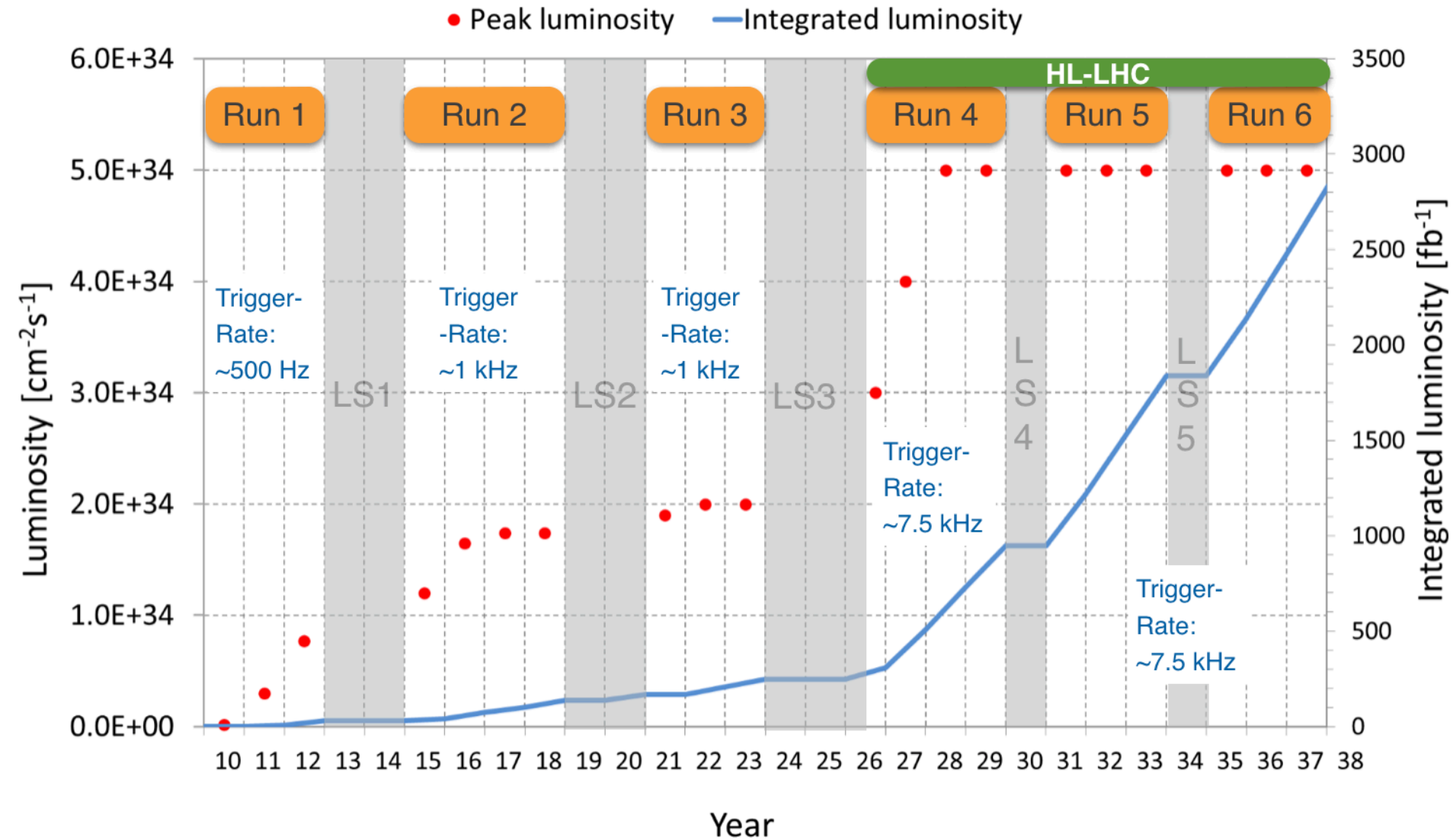**Deep neural network based algorithms perform the best**

## Networks can grow bigger, number of networks will increase

Network inferencing taking significant fraction of the final event processing time in CMS



Zoom ×10

DeepAK8

High top $p_T$

arXiv:1605.07678

# GROWING DATASET

Current: ~5 minutes per HL-LHC event



CMS Experiment at the LHC, CERN
Data recorded: 2016-Sep-08 08:30:28.497920 GMT
Run / Event / LS: 280327 / 55711771 / 67



**>5x**

Event complexity

**50x**

Total data
processing time

CPU seconds by Type

- Analysis
- HL-LHC MC
- LHC MC
- Non-Prompt Data
- Prompt Data

2027 estimate **CPU**:
~ 3.5 Million cores

Moore's Law continues …but Dennard Scaling fails

Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

Single threaded performance not improving
**Circa ~2005: "The Era of Multicore"**

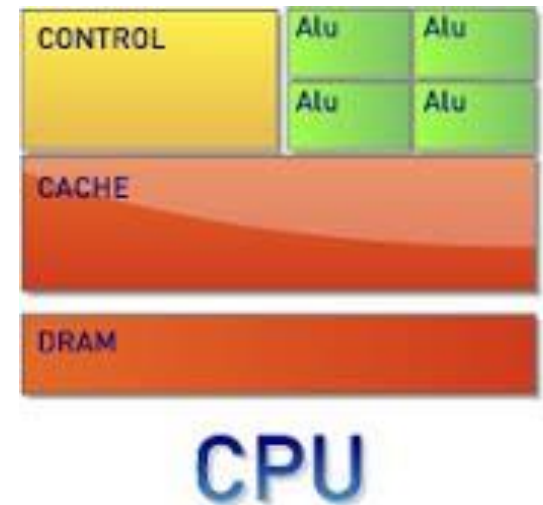Moore's Law continues …but Dennard Scaling fails

**We are not the only one facing the computing challenges faced with AI boom and data volume explosion**

Single threaded performance not improving

**Circa ~2005: "The Era of Multicore"**

→ **Today: Transition to the "Era of Specialization"?** (c.f. Doug Burger)

CMS week 2016    CERN June 20th - 24th

Remember that Facebook ask you (at least used to) to tag people when you upload a photo?

CMS week 2016   CERN June 20th - 24th

Runs image detection every time some one uploads a photo:
Neutral network inference

# CMS PARTY@2016



**CMS week 2016**   CERN June 20th - 24th

**300 million photos uploaded/day as of 2018.Nov**

# PLATFORM PROS & CONS:

| | | Perf/W | | |
|---|---|---|---|---|
| **More Flexible** | CPUs | 1X | Today's standard, most programmable, good for services changing rapidly | Conventional programming |
| | Manycore CPUs | 3X | Many simple cores (10s to 100s per chip), useful if software can be fine-grain parallel, difficult to maintain. | |
| | GPUs | 5-30X | Good for data parallelism by merged threads (SIMD), High memory bandwidth, power hungry | |
| | **FPGAs** | 5-30X | Most radical fully programmable option. Good for streaming/irregular parallelism. Power efficient but currently need to program in H/W languages. | Alternative programming |
| | Structured ASICS | 20-100X | Lower-NRE ASICs with lower performance/efficiency. Includes domain-specific (programmable) accelerators. | |
| **More Efficient** | Custom ASICs | > 100X | Highest efficiency. Highest NRE costs. Requires high volume. Good for functions in very widespread use that are stable for many years. | Can't change functionality |

Homogeneous — Specialized (vertical axis)

C/C++ — CUDA — Verilog — Verilog (right axis)

More Flexible

Homogeneous

Specialized

More Efficient

| | | Perf/W | |
|---|---|---|---|
| | CPUs | 1X | Today's standard, most programmable, good for services changing rapidly |
| | Manycore CPUs | 3X | Many simple cores (10s to 100s per chip), useful if software can be fine-grain parallel, difficult to maintain. |
| | GPUs | 5-30X | Good for data parallelism by merged threads (SIMD), High memory bandwidth, power hungry |
| | **FPGAs** | 5-30X | Most radical fully programmable option. Good for streaming/irregular parallelism. Power efficient but currently need to program in H/W languages. |
| | Structured ASICS | 20-100X | Lower-NRE ASICs with lower performance/efficiency. Includes domain-specific (programmable) accelerators. |
| | Custom ASICs | > 100X | Highest efficiency. Highest NRE costs. Requires high volume. Good for functions in very widespread use that are stable for many years. |

**Power bill**

Conventional programming — C/C++

Alternative programming — CUDA / Verilog

Can't change functionality — Verilog

**Software/Electrical engineer Salaries**

Specialized co-processor hardware for machine learning inference

ASIC?

A11 Bionic neural engine

A11

FPGA

Stratix 10
1SG280LN3F43E3VGS1
US

Microsoft

FPGA

Catapult/Brainwave

intel Microsoft

Delivering FPGA Partner Solutions on AWS
via AWS Marketplace

Customers

AWS Marketplace

FPGA

Amazon Machine Image (AMI)

Amazon FPGA Image (AFI)

AFI is secured, encrypted, dynamically loaded into the FPGA - can't be copied or downloaded

GA
tplace

Google
Tensor Processing Unit

ASIC

**INTEL® FPGA ACCELERATION HUB**

The Intel® Xeon® Acceleration Stack for FPGAs is a robust framework enabling data center applications to leverage an FPGA's potential to increase

Computationally intensive: iterative algorithms such as track reconstruction

**Option 1**

**re-write physics algorithms for new hardware**

Language: OpenCL, OpenMP,TBB, HLS, …?

Hardware: FPGA, GPU

**Option 2**

**re-cast physics problem as a machine learning problem**

Language: C++, Python (TensorFlow, PyTorch,…)

Hardware: FPGA, GPU, ASIC

**Example: tracking@HL-LHC:**

**Option 1: Parallelized and Vectorized Tracking Using Kalman Filters**

**Option 2: Recent work on tracking using Graph Networks**

# POSSIBLE SOLUTIONS FOR US

**Option 1**

**re-write physics algorithms for new hardware**

Language: OpenCL, OpenMP,TBB, HLS, …?

Hardware: FPGA, GPU

**Option 2**

**re-cast physics problem as a machine learning problem**

Language: C++, Python (TensorFlow, PyTorch,…)

Hardware: FPGA, GPU, ASIC

**Advantage of option 2: recasting problem as machine learning problems (computing wise)**

- Algorithms can universally be expressed as simple matrix multiplications computations
- Intrinsically parallelizable
- Follow industry trends in developing co-processors optimized for ML and speed the up the inference(sub-event level reconstruction such as tracking)

# Proof of concept: Particle physics computing with Brainwave

Will explain this later with
pretty pictures
Picked this because of its
mature eco system

**1μs**      **1ms**      **1s**

**LHC L1 Trigger
(pipelined)**
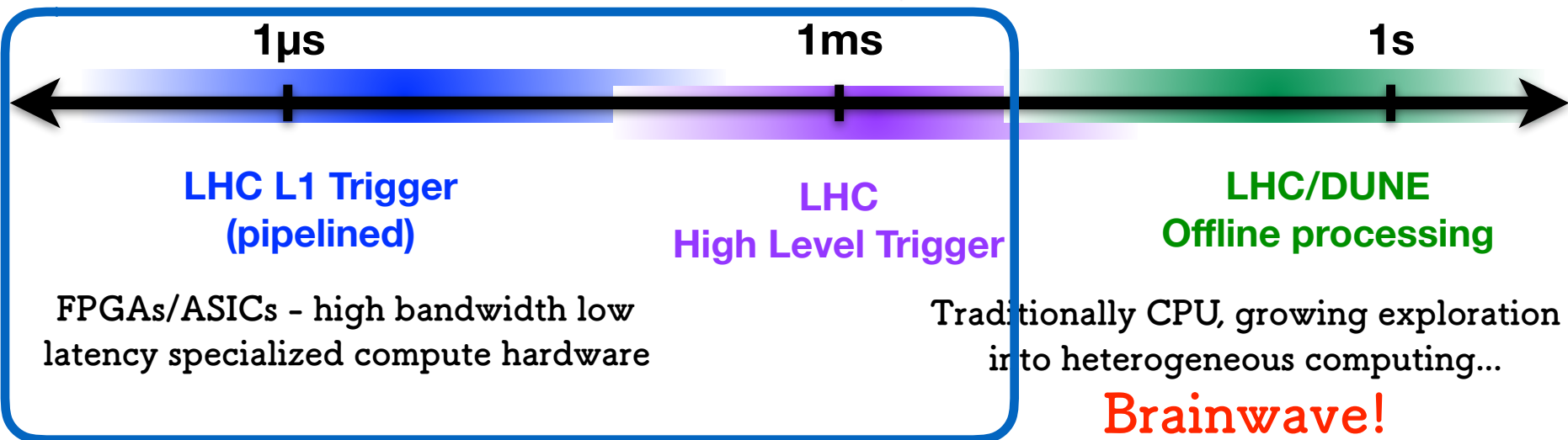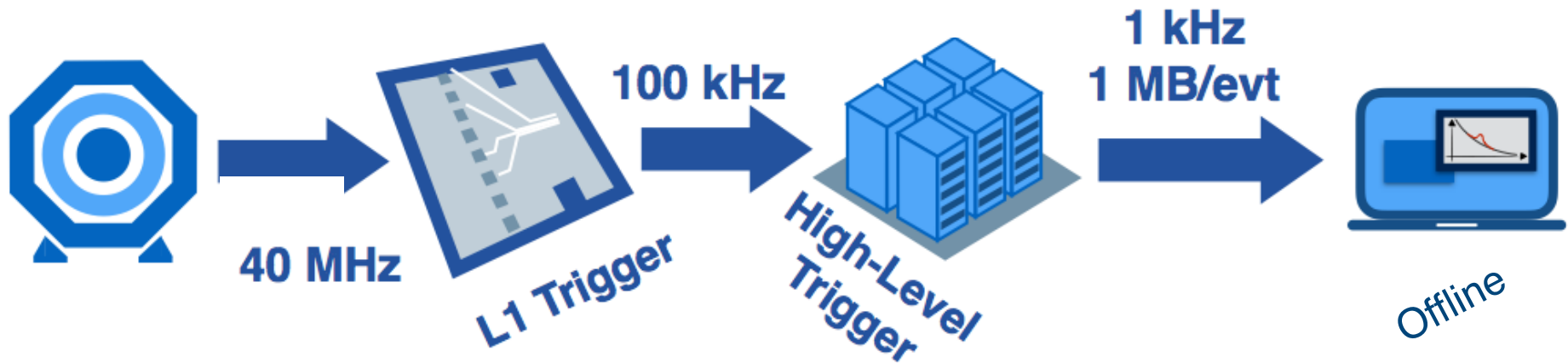
**LHC
High Level Trigger**

**LHC/DUNE
Offline processing**

FPGAs/ASICs - high bandwidth
low latency specialized
compute hardware

Traditionally CPU, growing
exploration into heterogeneous
computing...

**1μs**      **1ms**      **1s**

**LHC L1 Trigger (pipelined)**

**LHC High Level Trigger**

**LHC/DUNE Offline processing**

FPGAs/ASICs – high bandwidth low latency specialized compute hardware

Traditionally CPU, growing exploration into heterogeneous computing...

**Brainwave!**

**Two parallel talks this afternoon for L1/HLT applications:**
"DNN based algorithm for CMS Level-1 muon reconstruction" by Jia Low.
"Deep Machine Learning on FPGAs for L1 trigger and Data Acquisition" by Dylan Rankin

Even if co-processors are 100x faster, is it feasible to have every T1,T2,T3 computing farm buy specialized hardware?

**No, but…**
**Interesting possibility for the HLT farm…**



**Offline solution**: co-processors as a service

Hardware as a Service
Network Acceleration
Compute Acceleration

Stratix FPGA

40G Ethernet

NIC

CPU

PCIe Gen 3

PCIe Gen 3

**For more on MS catapult: see talk by A. Putnam**
**https://www.dropbox.com/s/rvd06vp5ogguqxe/Catapult_2018_Fermilab_Public.pdf**

# PROOF OF CONCEPT: SONIC

**S**ervices for **O**ptimized **N**etwork **I**nference on **C**o-processors

*PRELIMINARY RESULTS!*

(work in progress)

**FPGA-accelerated machine learning inference as a solution for particle physics computing challenges**

Javier Duarte, Burt Holzman, Ben Kreis, Mia Liu, Kevin Pedro, N.T., Aris Tsaris (FNAL)
Phil Harris, Dylan Rankin (MIT)

+ Doug Burger, Eric Chung, Andrew Putnam (MS research), Ted Way,MS, David Lee (MS Azure)

**Question:**
How do we integrate heterogeneous computing resources into the physics event data processing model?

# ACCESSING HETEROGENEOUS RESOURCES

Implemented New CMSSW feature called ExternalWork:

○ Asynchronous task-based processing
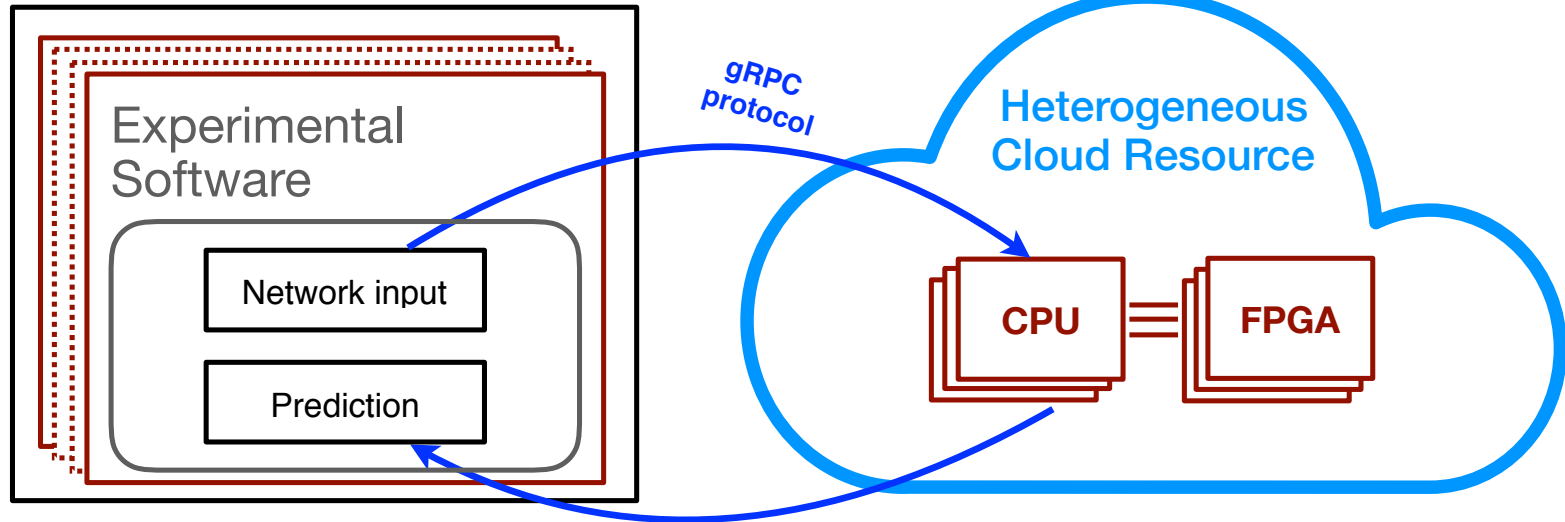


○ **Non-blocking**: schedule other tasks while waiting for external processing
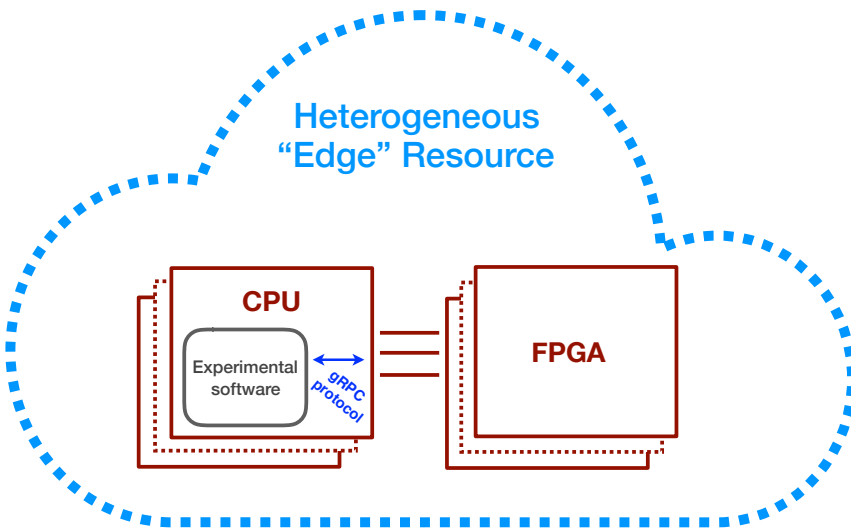
Can be used with GPUs, FPGAs, cloud, …

➢**Now demonstrated to work with Microsoft Brainwave!**

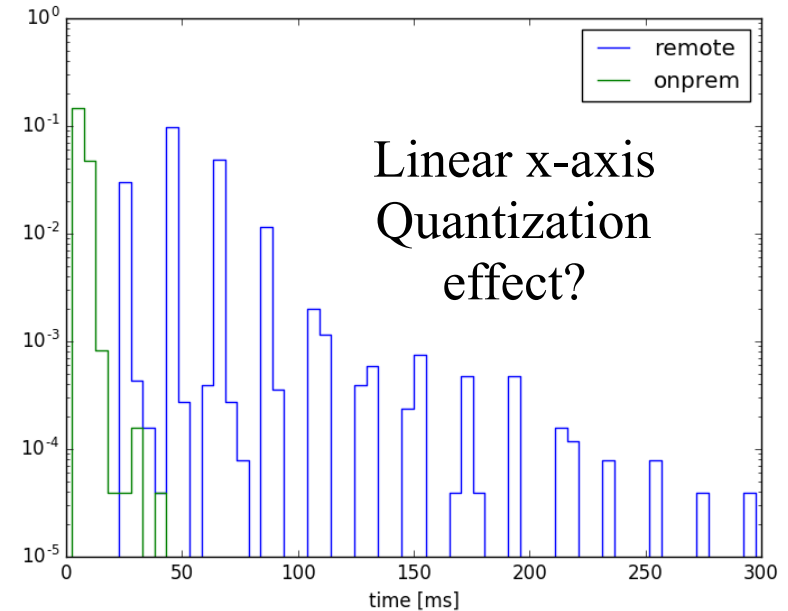More details on external work module: Kevin Pedro's talk at CHEP

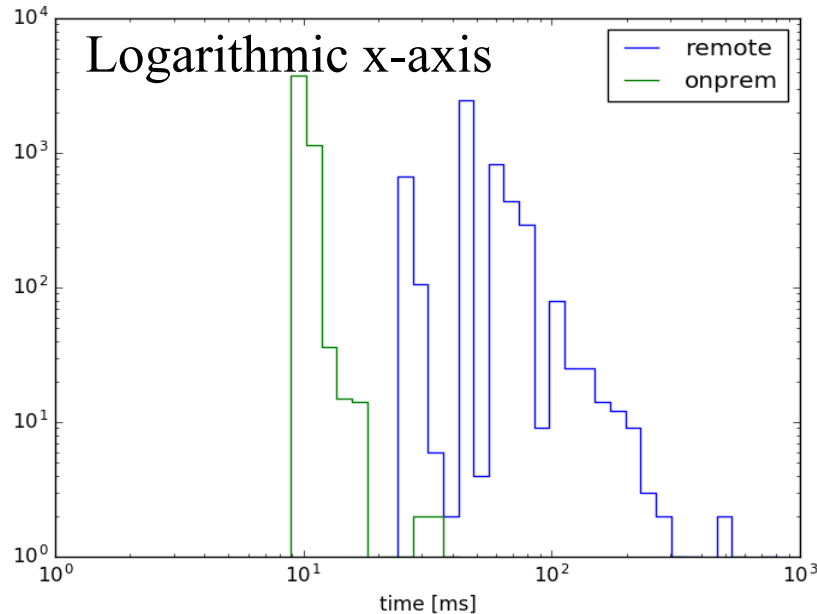**Datacenter (CPU farm)**



- Cloud service has latency

- Run CMSSW on Azure cloud machine → simulate local installation of FPGAs ("on-prem" or "edge")
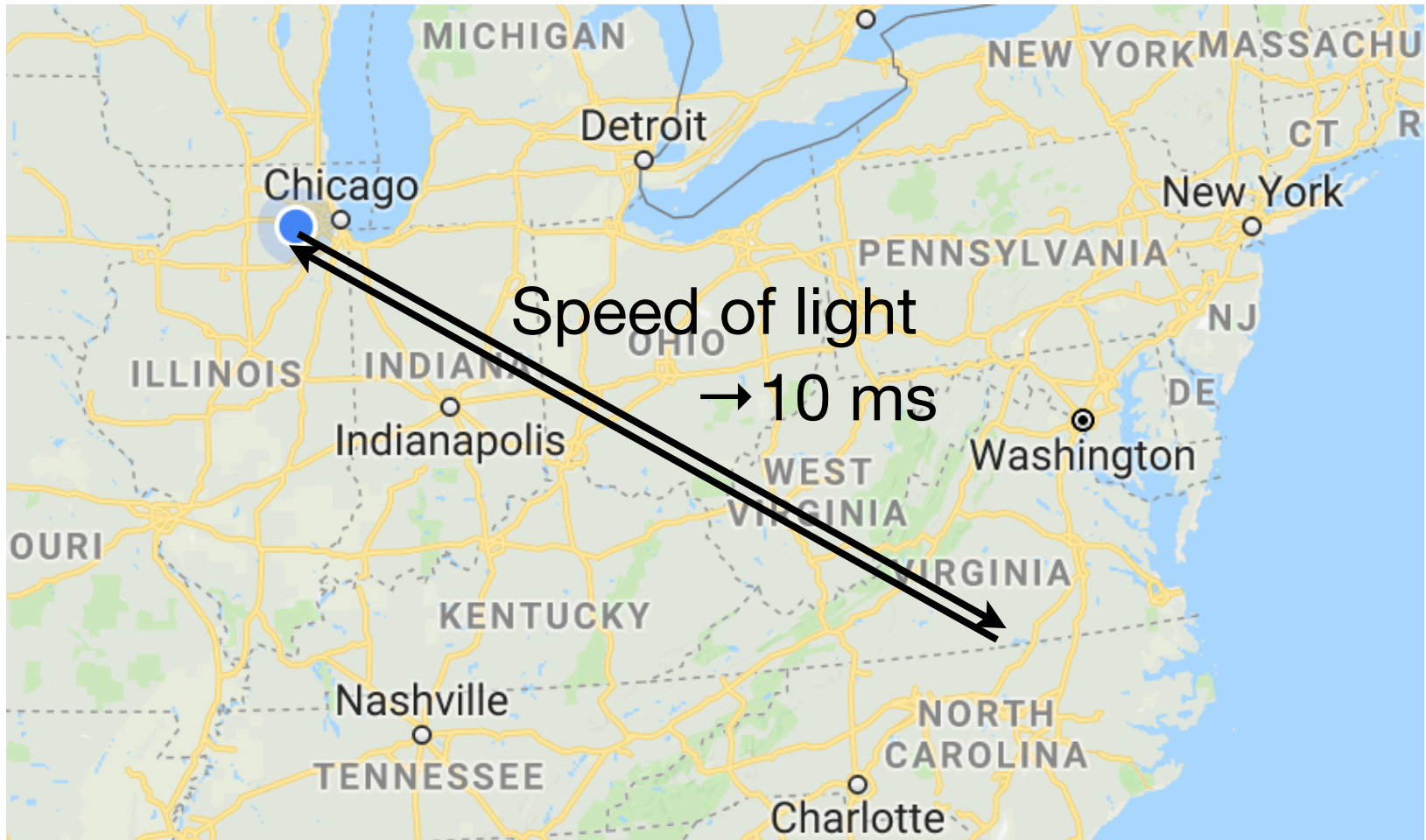
- Provides test of "HLT-like" performance

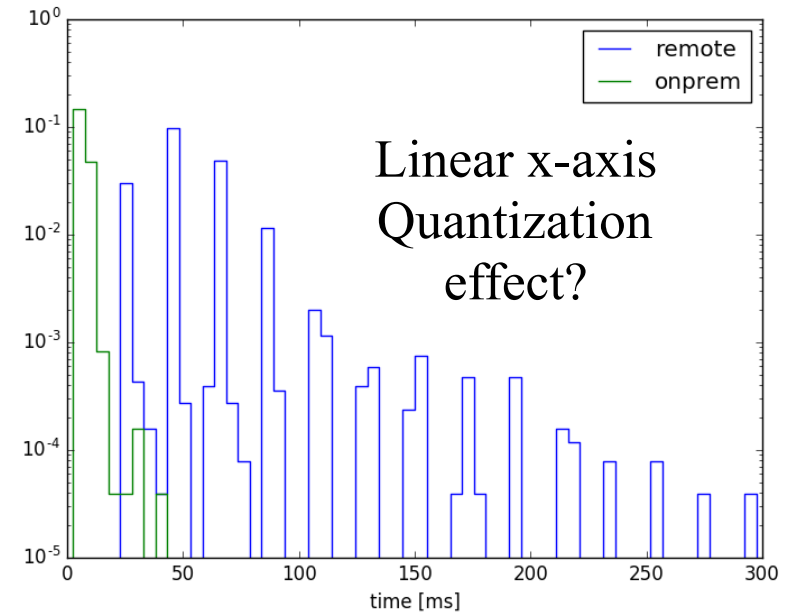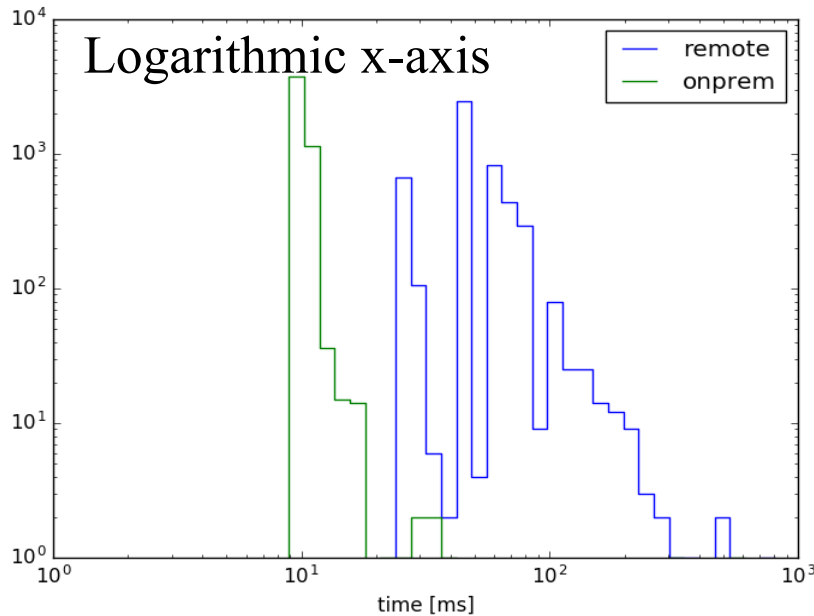Good performance in initial tests
- "remote": cmslpc @ FNAL to Azure (VA),    ‹time› = 56 ms
- "onprem": run CMSSW on Azure VM,    ‹time› = 10 ms
  (~2 ms on FPGA, rest is classifying and I/O)

Speed of light → 10 ms

With network switches?  May be about right :)

Logarithmic x-axis — remote, onprem

Linear x-axis
Quantization effect?

Good performance in initial tests

○ "remote": cmslpc @ FNAL to Azure (VA),   ‹time› = 56 ms

○ "onprem": run CMSSW on Azure VM,        ‹time› = 10 ms
   (~2 ms on FPGA, rest is classifying and I/O)

| Type | Hardware | Mean inference time | Setup |
|------|----------|---------------------|-------|
| CPU | Xeon 2.6 GHz, 1 core | 1.75 seconds | `CMSSW, TF v1.06` |
| CPU | i7 3.6 GHz, 1 core | 500 ms | standalone python, TF `v1.10` |
| CPU | i7 3.6 GHz, 8 core | 200 ms | standalone python, TF `v1.10` |

The chart legend:
- azure resnet gpu
- resnet gpu
- resnet gpu train

*NVidia GTX 1080 Ti*

"onprem FPGA": 10 ms

Brainwave ResNet50 on GPU

Official ResNet50 in tensorflow

Super optimized ResNet50

Not so straightforward to compare against other hardware, the whole chain matters: pipelined inputs, IO bandwidth (PCIe), special instruction sets, etc. General findings:

**GPUs:** O(~100 ms), for batch-1 input

To explore: Google TPUs, AWS/Xilinx FPGAs, Intel/Altera FPGAs

Exploring the use of FPGA co-processors (MS Brainwave) for ML acceleration as an "off-the-shelf" computing paradigm for particle physics

- Deploying cloud accelerators as a service **fits the particle physics computing model in a non-disruptive way**

    - For large computing tasks (Resnet-50), **there is a ~(4/10/100)x benefit over CPU-only computations**

    - Could be used for neutrino experiments **~today!**

- "Edge" compute option as an HLT solution?

*Outlook and further studies*

To explore: Google TPUs, AWS/Xilinx FPGAs, Intel/Altera FPGAs
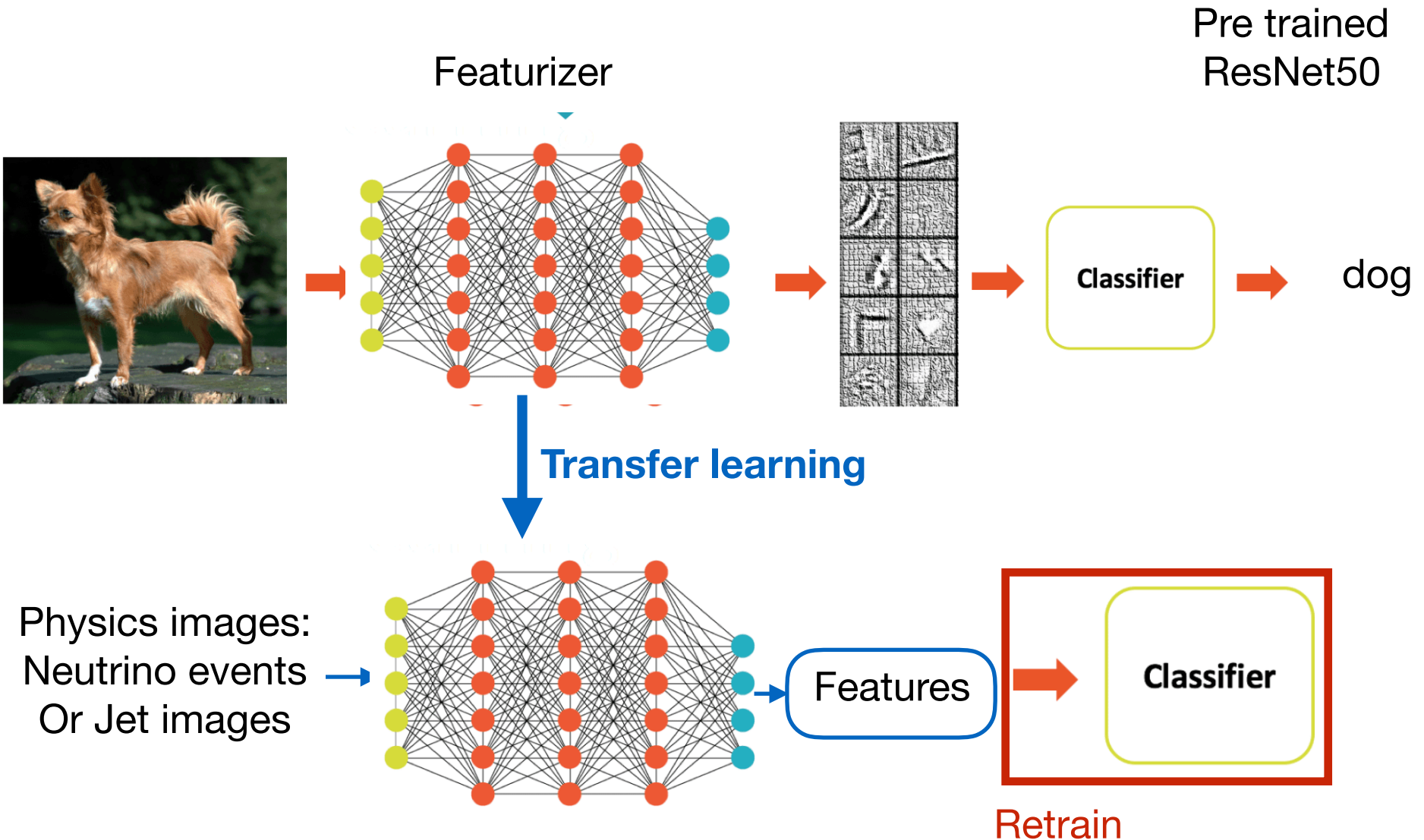
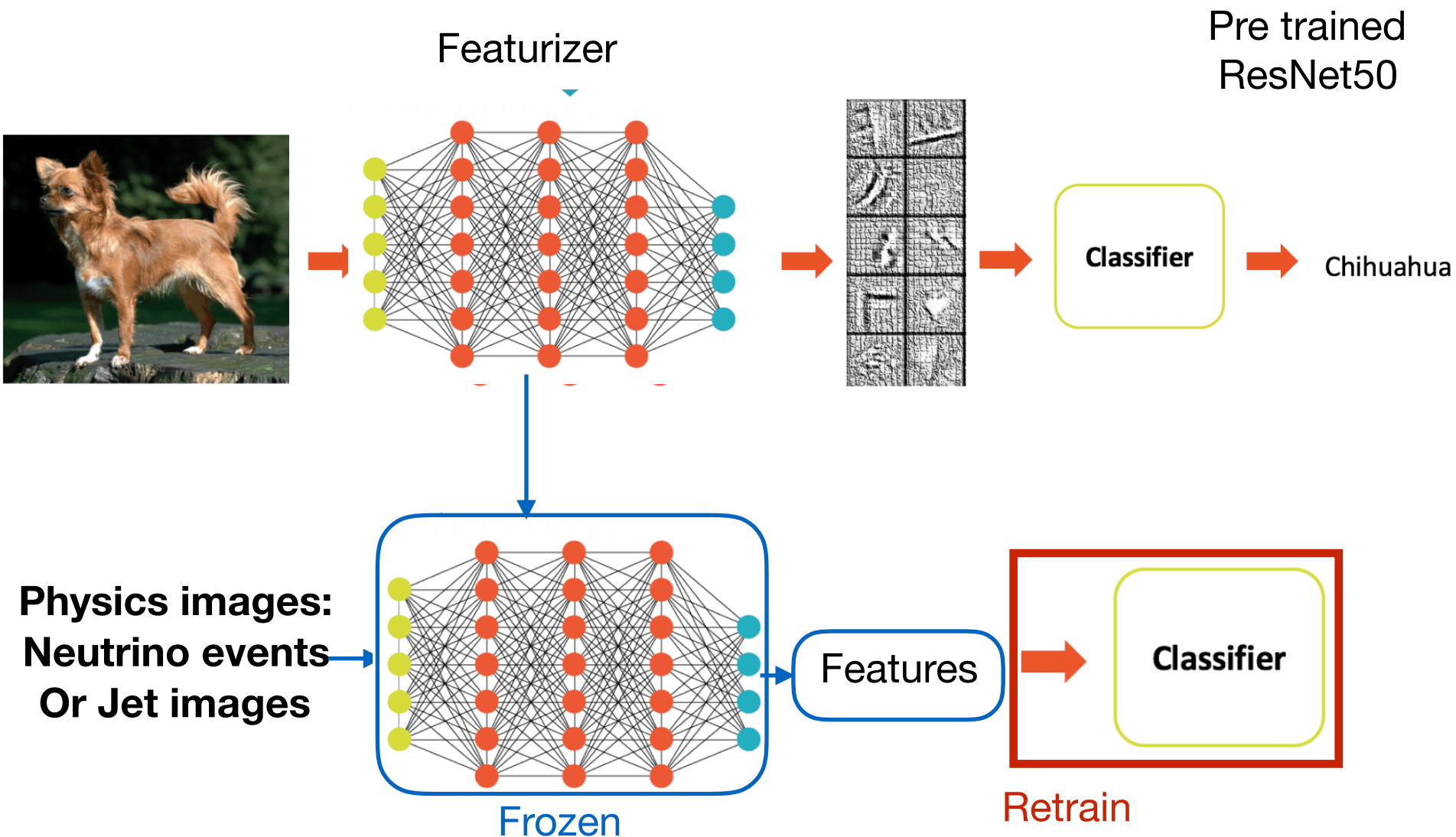    Important to benchmark different platforms to understand our options/projections.

Need to understand scaling (not too worried about this)
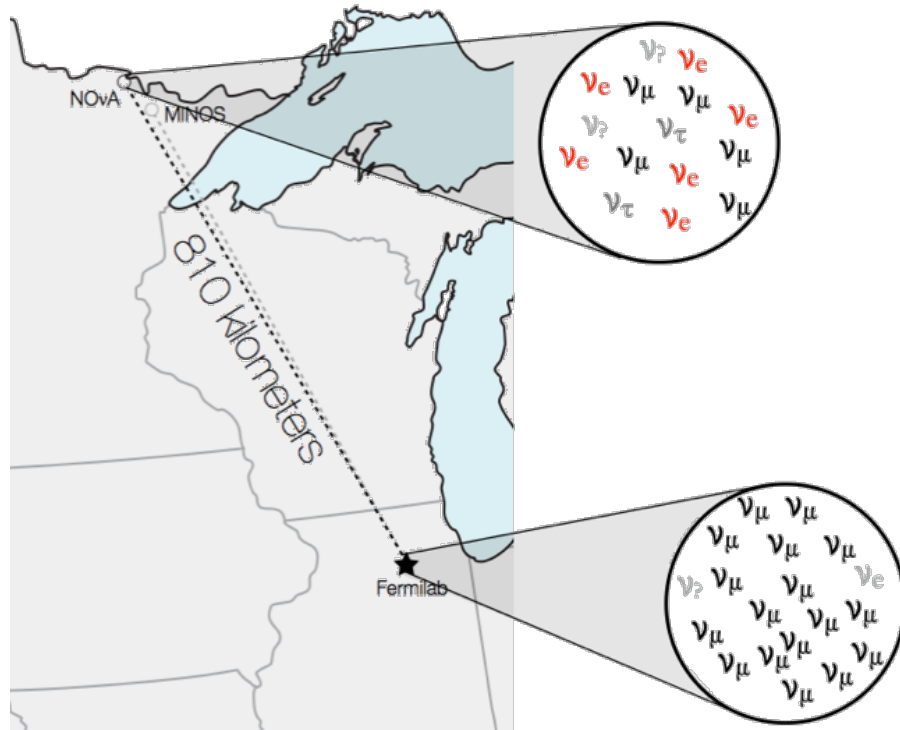
What's the costing model?

Relies on continued development of ML algorithms for difficult physics problems (simulation/reconstruction)

# Physics case: Event classification in NOvA

Featurizer

Pre trained
ResNet50

Classifier

dog

**Transfer learning**

Physics images:
Neutrino events
Or Jet images

Features

Classifier

Retrain

Featurizer

Pre trained
ResNet50

Classifier → Chihuahua

Physics images:
Neutrino events
Or Jet images

Features

Classifier

Frozen

Retrain

New feature: fine-tune the weights in featurizer too! Will be included in final results
Other models became available recently: VGG etc

# THE NOvA EXPERIMENT:

NuMI: Neutrinos at the Main Injector

Long-baseline (anti-)neutrino oscillation experiment

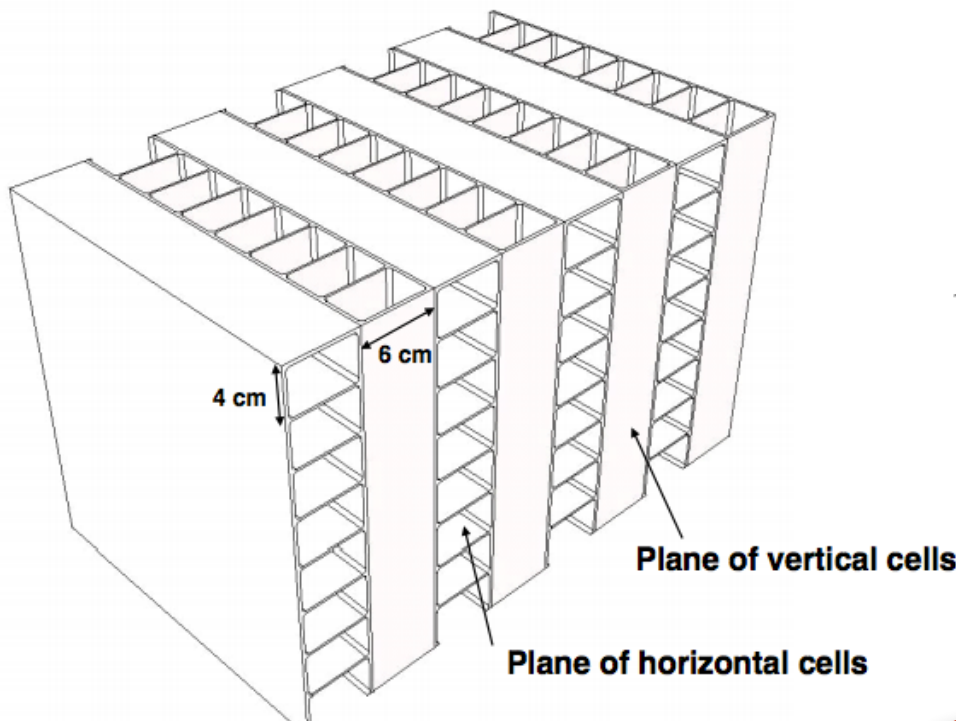Two functionally identical detectors, optimized for $\nu_e$ identification

**Primary goal: measurement of neutrino oscillations via $\nu\mu\rightarrow\nu e$**
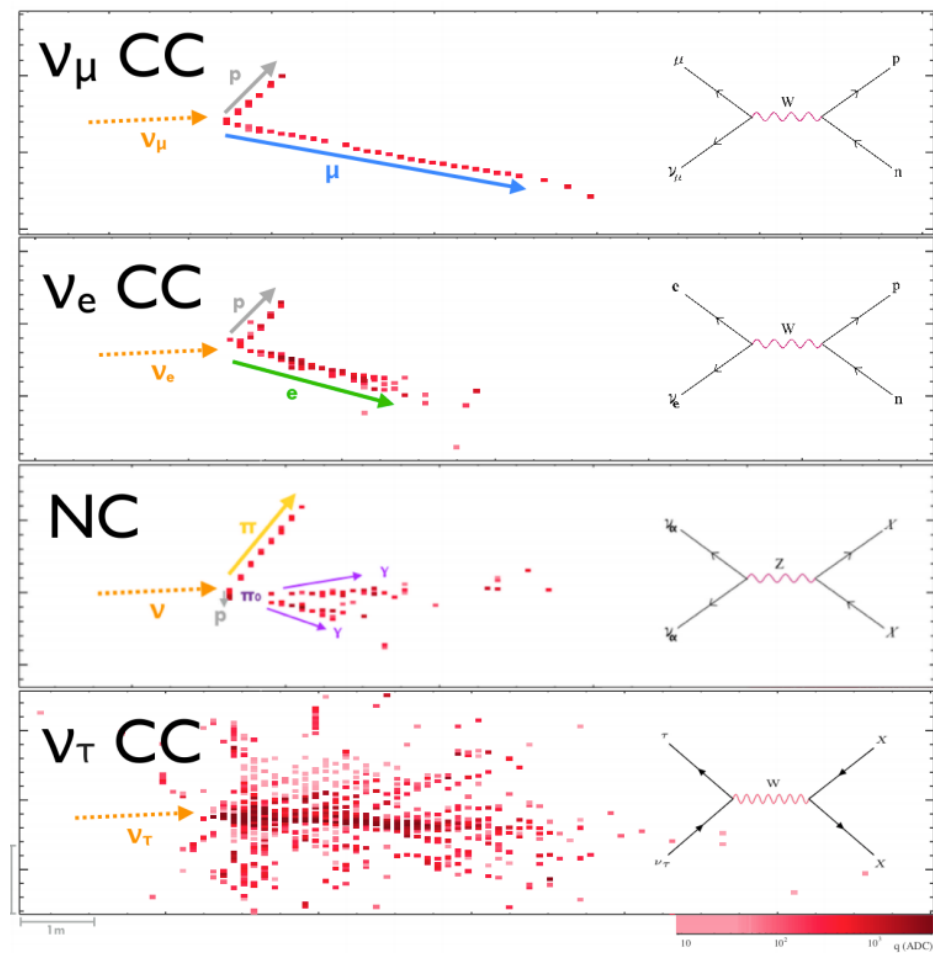
Other goals include: Searches for sterile neutrinos Neutrino cross sections Supernova neutrinos Cosmic ray physics

6 cm

4 cm

**Plane of vertical cells**

**Plane of horizontal cells**

$\nu_\mu$ CC

$\nu_e$ CC

NC

$\nu_\tau$ CC

**3D reconstruction**

Far Detector, on surface
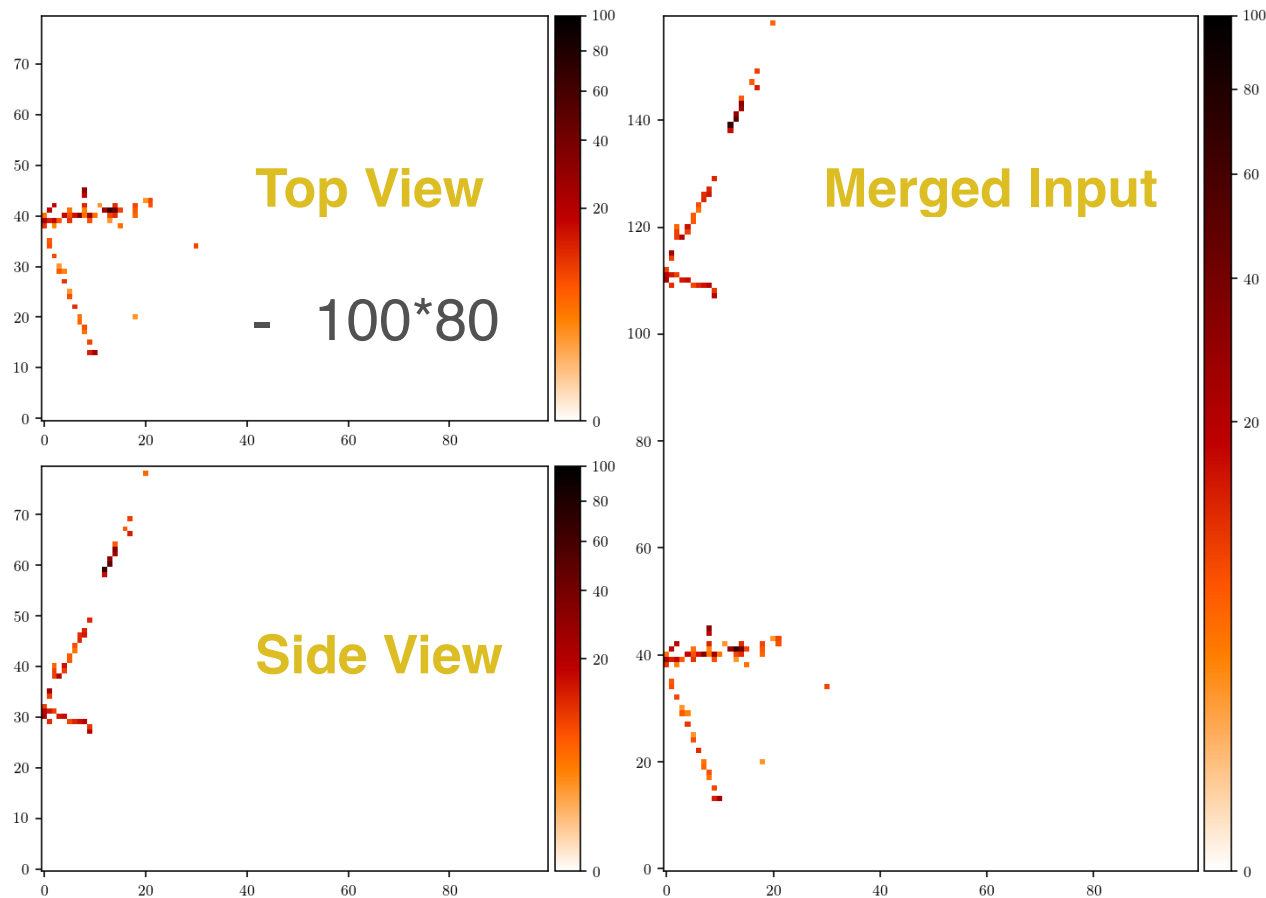**14 kTons
896 readout planes**

Training Setup
- training/testing: 500K/150K
- Pre-trained ResNet-50 model on image net.

**5 labels:**
- Muon neutrino
- Electron neutrino
- Tau neutrino
- Neutral Current
- Cosmic



**Top View**
- 100*80

**Side View**

**Merged Input**

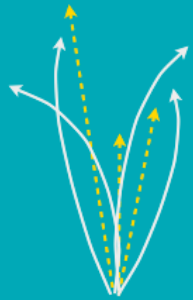- Merged image scaled to resolution of 224*224 using Bilinear Interpolation from TF, to be fed into ResNet50

Electron Neutrino — Muon Neutrino — Tau Neutrino — Cosmic — Neutral Current (Top View / Side View)

Events identified with more that 0.9 probability by the ResNet-50 network. Color represents energy deposit
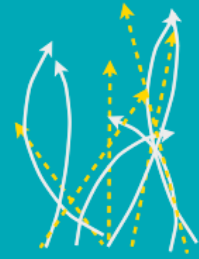
# Physics cases: jet substructure

*u,d* or *s* jet

*c* or *b* jet

gluon jet

pileup jet

*W* or *Z* jet

Higgs jet

top jet

?

# TRAINING STATUS
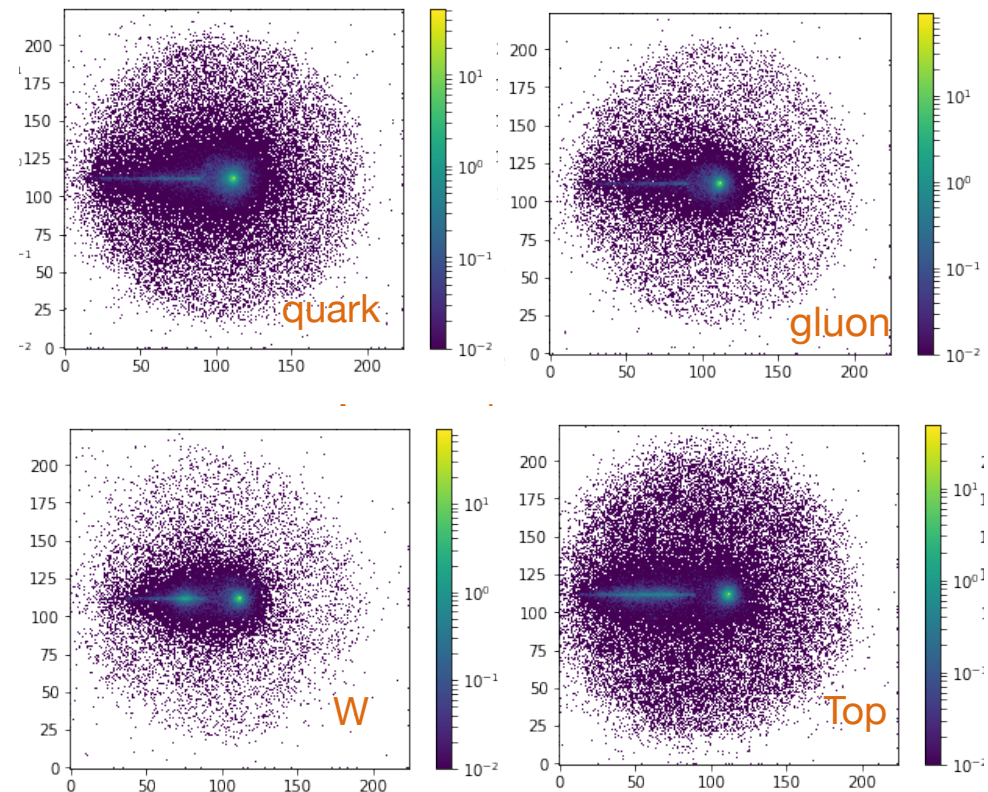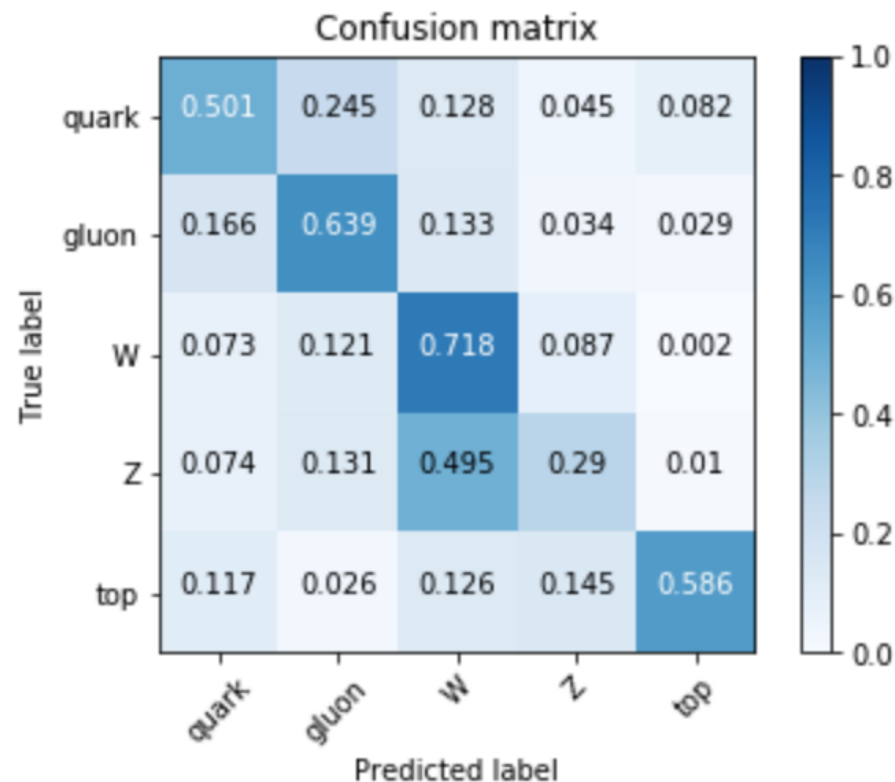
Generator level AK8 jets: quark/gluon/W/Z/top, density map of the pt of jet constituents.



Averaged over 1000 images
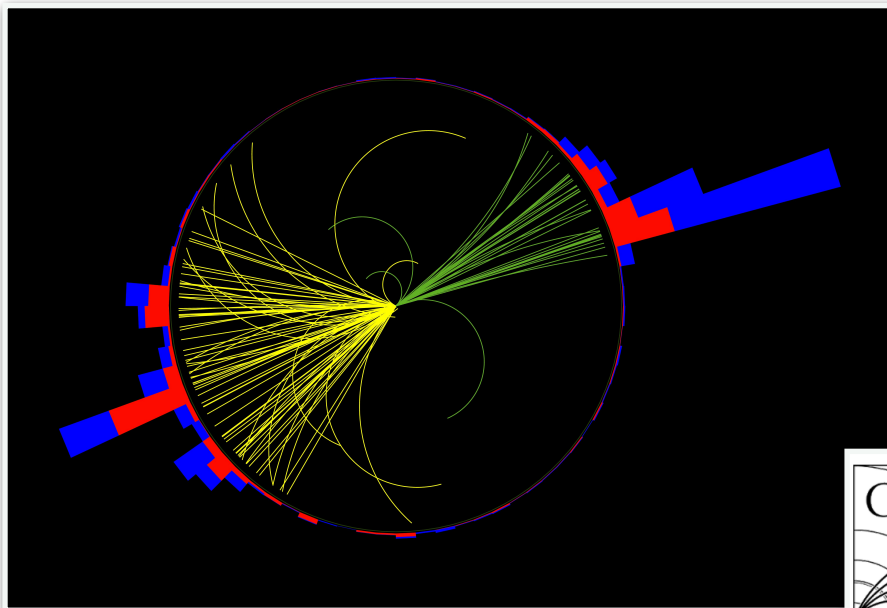
# Outlook and next steps

# TAKEAWAYS

Computing challenges in big data: More complex detectors and sophisticated algorithms, large datasets.

We follow the industry trend in exploring specialized hardware (co-processers) as ML acceleration options.
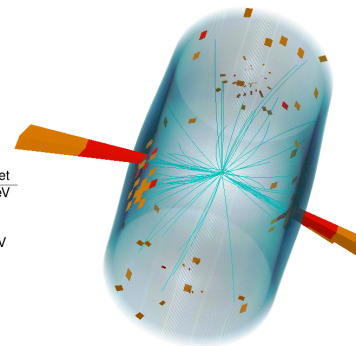
Started with Microsoft brainwave, **demonstrated FPGAs are a promising option to accelerate neural network inference:**

o *Can achieve (at least) order of magnitude improvement over CPU*

o *Better fit for CMS event-level computing model (vs. GPUs which require batching for efficiency)*

o *Physics cases: Nova event classification& jet substructure using ResNet50 on brainwave.*

**Proof of concept, more studies to follow**

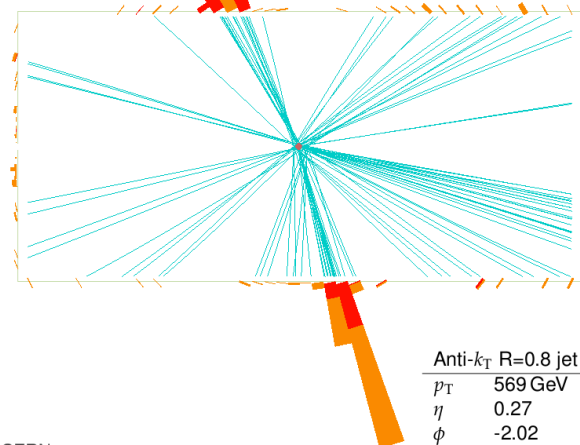# BACKUP

**Candidate qW event**
Dijet mass: 5.1 TeV

Anti-$k_T$ R=0.8 jet
| | |
|---|---|
| $p_T$ | 2406 GeV |
| $\eta$ | 0.66 |
| $\phi$ | 2.51 |
| $M_{SD}$ | 29.1 GeV |
| $\tau_{21}$ | 0.50 |

Anti-$k_T$ R=0.8 jet
| | |
|---|---|
| $p_T$ | 2298 GeV |
| $\eta$ | -0.17 |
| $\phi$ | -0.63 |
| $M_{SD}$ | 81.6 GeV |
| $\tau_{21}$ | 0.29 |

**Candidate Z jet**

Anti-$k_T$ R=0.8 jet
| | |
|---|---|
| $p_T$ | 2.1 TeV |
| $\eta$ | -0.32 |
| $\phi$ | 0.63 |
| $M_{SD}$ | 96.6 |
| $\tau_2/\tau_1$ | 0.34 |

CMS Experiment at LHC, CERN
Data recorded: Sat Oct 22 00:05:32 2016 CEST
Run/Event: 283820 / 450110972
Lumi section: 263

Anti-$k_T$ R=0.8 jet
| | |
|---|---|
| $p_T$ | 618 GeV |
| $\eta$ | -0.53 |
| $\phi$ | 1.18 |
| $M_{SD}$ | 81.3 GeV |
| $\tau_{21}$ | 0.29 |

**Candidate WW event**
Dijet mass: 1.3 TeV

Anti-$k_T$ R=0.8 jet
| | |
|---|---|
| $p_T$ | 569 GeV |
| $\eta$ | 0.27 |
| $\phi$ | -2.02 |
| $M_{SD}$ | 80.2 GeV |
| $\tau_{21}$ | 0.32 |

CMS Experiment at LHC, CERN
Data recorded: Fri Aug 19 02:26:23 2016 CEST
Run/Event: 279024 / 602168401
Lumi section: 376

# FURTHER NEXT STEPS

## Model customization

Whenever we talk about this — people are generally positive but always ask when we can put our own networks on the FPGAs

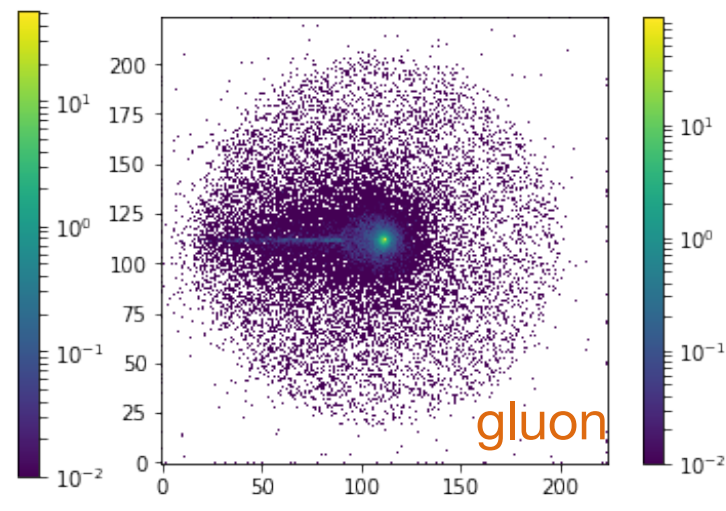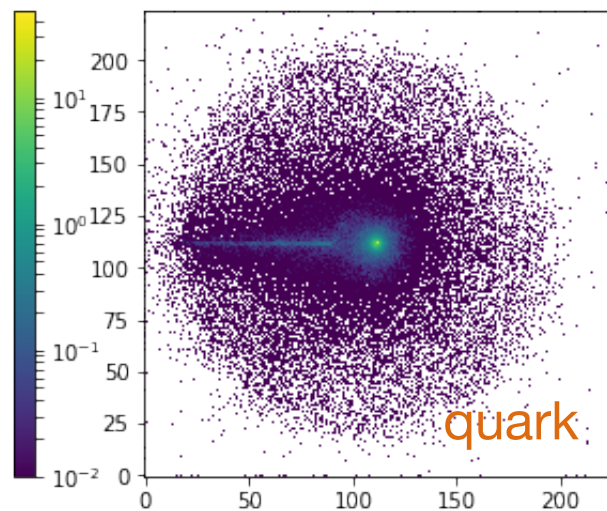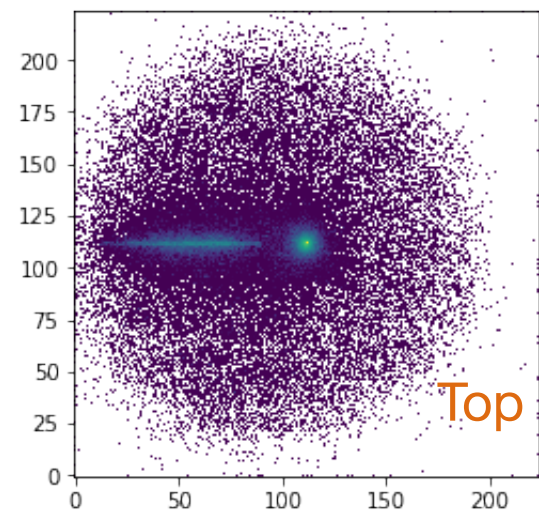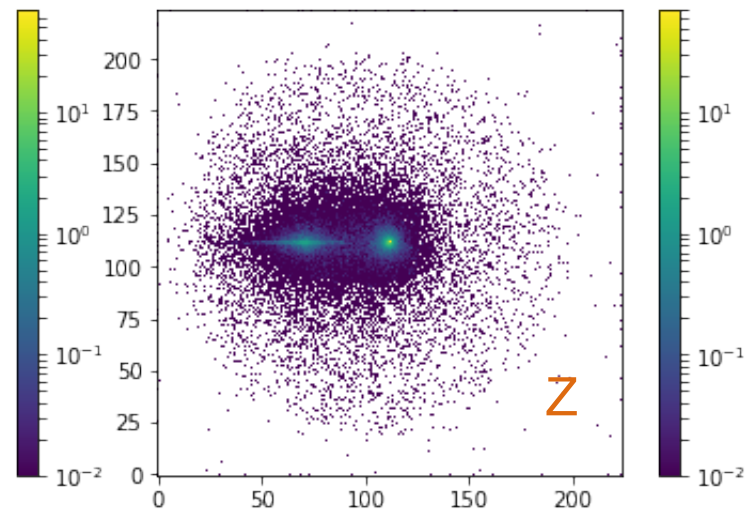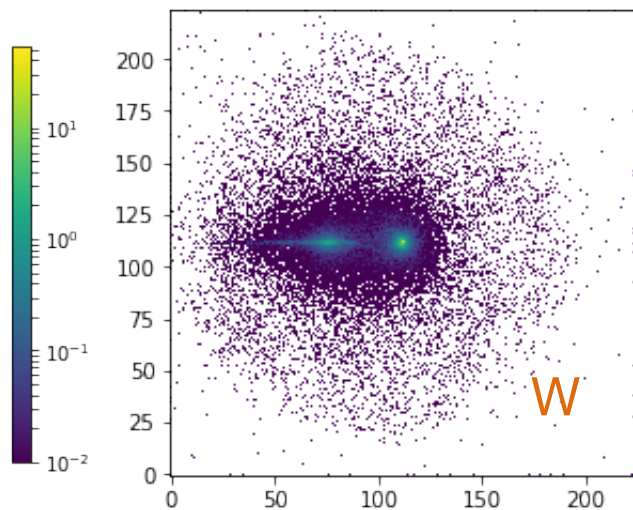Is it something we can work with you on?  Not just CNNs, but Graph NNs, LSTMs, etc…

## On-prem HLT-like demonstration
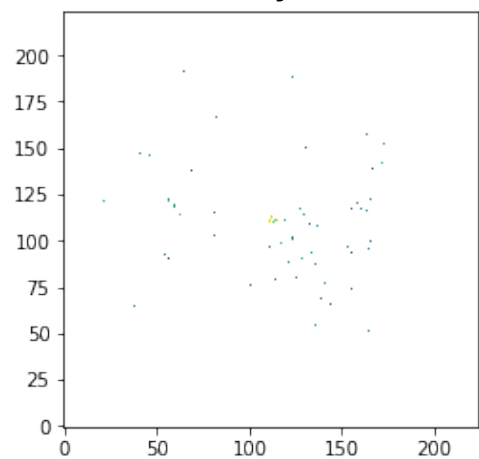
An "edge" offering has been brought up a few times — this is something, with necessary infrastructure, we're interested in pursuing if possible as a demonstration of the trigger (on-prem, real-time) capabilities

## Scaling up

We should try to demonstrate running on N (>>1) CPUs  and M (>1) services to understand how to scale services.  This will give us an idea of cost scaling as well.

1 W jet

W

Z

Top

quark

gluon

Averaged over 1000 images

# THE CMS SOFTWARE

**CMSSW:**

- Hosted on [GitHub](GitHub)

- ~6 million lines of code

- Handles simulation, raw data processing, reconstruction, analysis

**Event-based processing model:**

- Load event data into memory

- Numerous modules process parts of event, output new products

**Parallelism:**

- Multiple events in flight → *streams*

- Multiple modules running simultaneously → *threads*

  - Task-based multithreading using Intel Thread Building Blocks

**WLCG: Worldwide LHC Computing Grid**

- Network of computing clusters at labs, universities, etc.
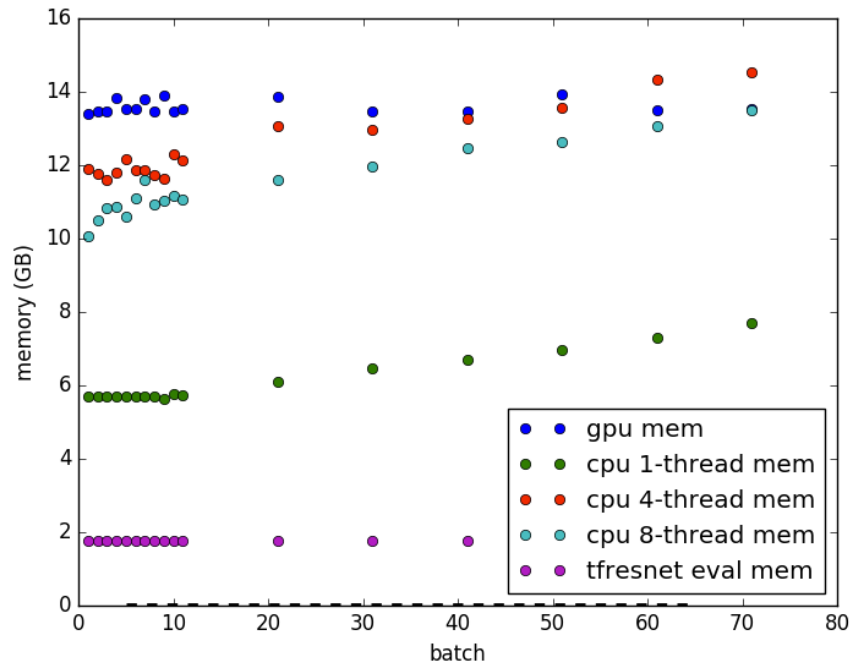
- Mostly commodity hardware

# SONIC IN CMSSW

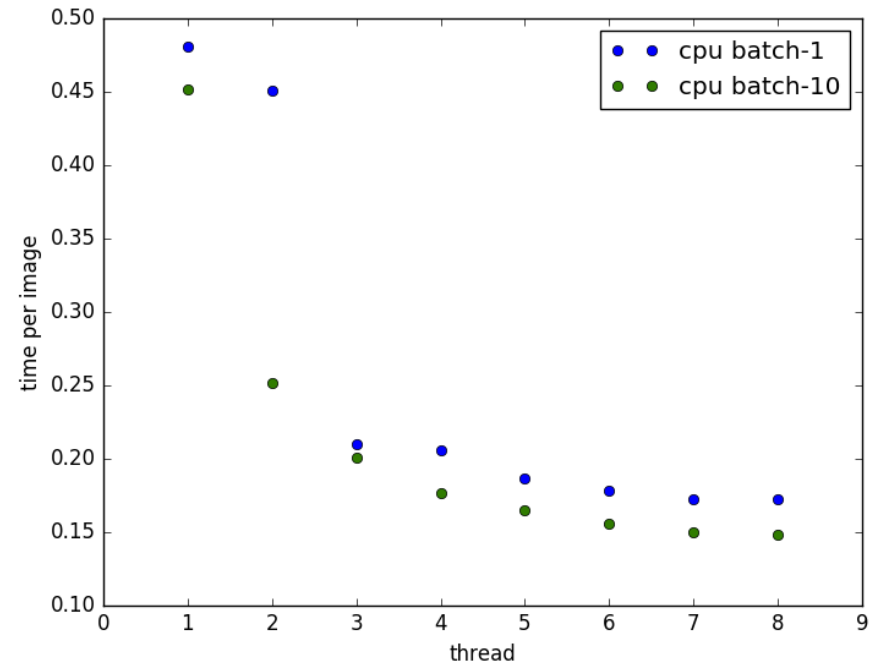**S**ervices for **O**ptimized **N**etwork **I**nference on **C**oprocessors

Demonstration in Microsoft Brainwave:

- ○ Create "image" from jet constituents, process with ResNet50

  - ▪ Much larger than custom HEP networks (so far)

- ○ Send to Microsoft Brainwave FPGA using gRPC w/ TensorFlow (protobuf)

- ○ FPGA processes one image at a time → no batching needed to be efficient

- ○ Use ExternalWork mechanism

  - ▪ gRPC C++ API lacks a callback interface (currently)

    - ➤ wait for gRPC return in lightweight std::thread
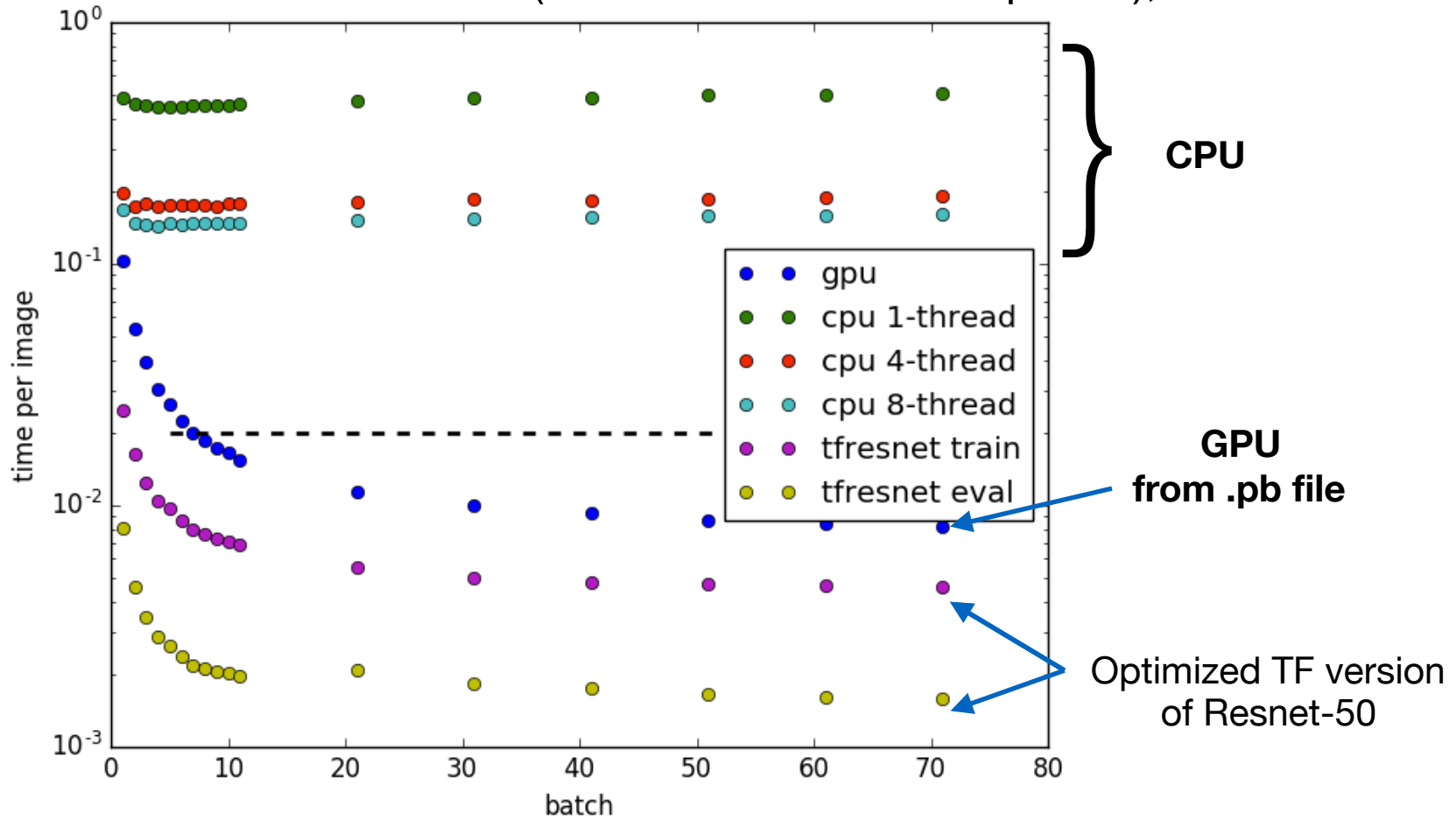
[SonicCMS](#) repository on GitHub

Memory usage in
different configurations

Latency versus number of threads

**Benchmark Nvidia GTX 1080, Intel i7 3.6 GHz**
Pure inference time (load time is 5 min for .pb file), TF v1.10



**CPU**

**GPU from .pb file**

gpu
cpu 1-thread
cpu 4-thread
cpu 8-thread
tfresnet train
tfresnet eval

Optimized TF version of Resnet-50

Full enqueuing with random inputs,
Large memory usage (12 Gb) with .pb input

**Benchmark Nvidia GTX 1080, Intel i7 3.6 GHz**
Pure inference time (load time is 5 min for .pb file), TF v1.10



CPU comparison:

**Intel i7 3.6 GHz (8 core, TF v1.10) ~ 180 ms**

**Intel i7 3.6 GHz (1 core, TF v1.10) ~ 500 ms**

**Intel i7 3.6 GHz (1 core, TF v1.06) ~ 1.2 s**

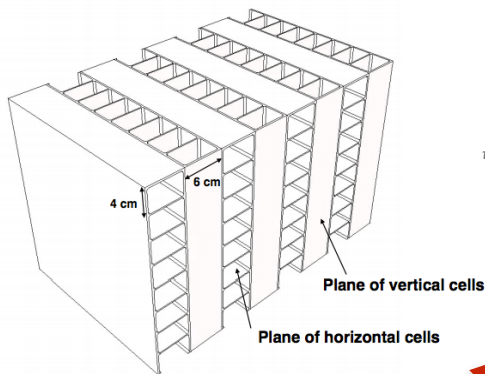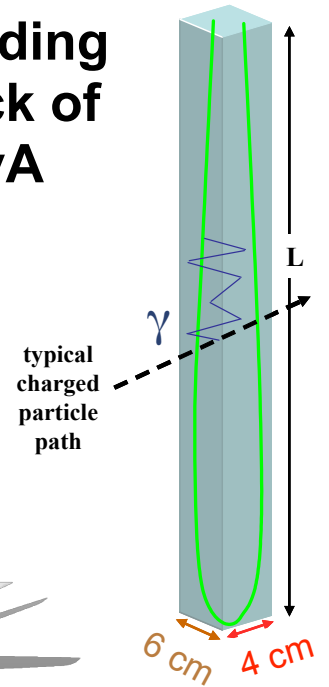**Intel Xeon 2.6 GHz (1 core, TF v1.06) ~ 1.75 s**
[what we are running]

Full enqueuing with random inputs,
Large memory usage (12 Gb) with .pb input

# NOvA DETECTORS

» Highly segmented low Z tracking calorimeter.

» Cells are filled wit liquid scintillator

  » wave shifting fiber readout.

» 65% active by volume

» Detection with avalanche photo diodes.

» Alternating X/Y planar geometry: **3D reconstruction**

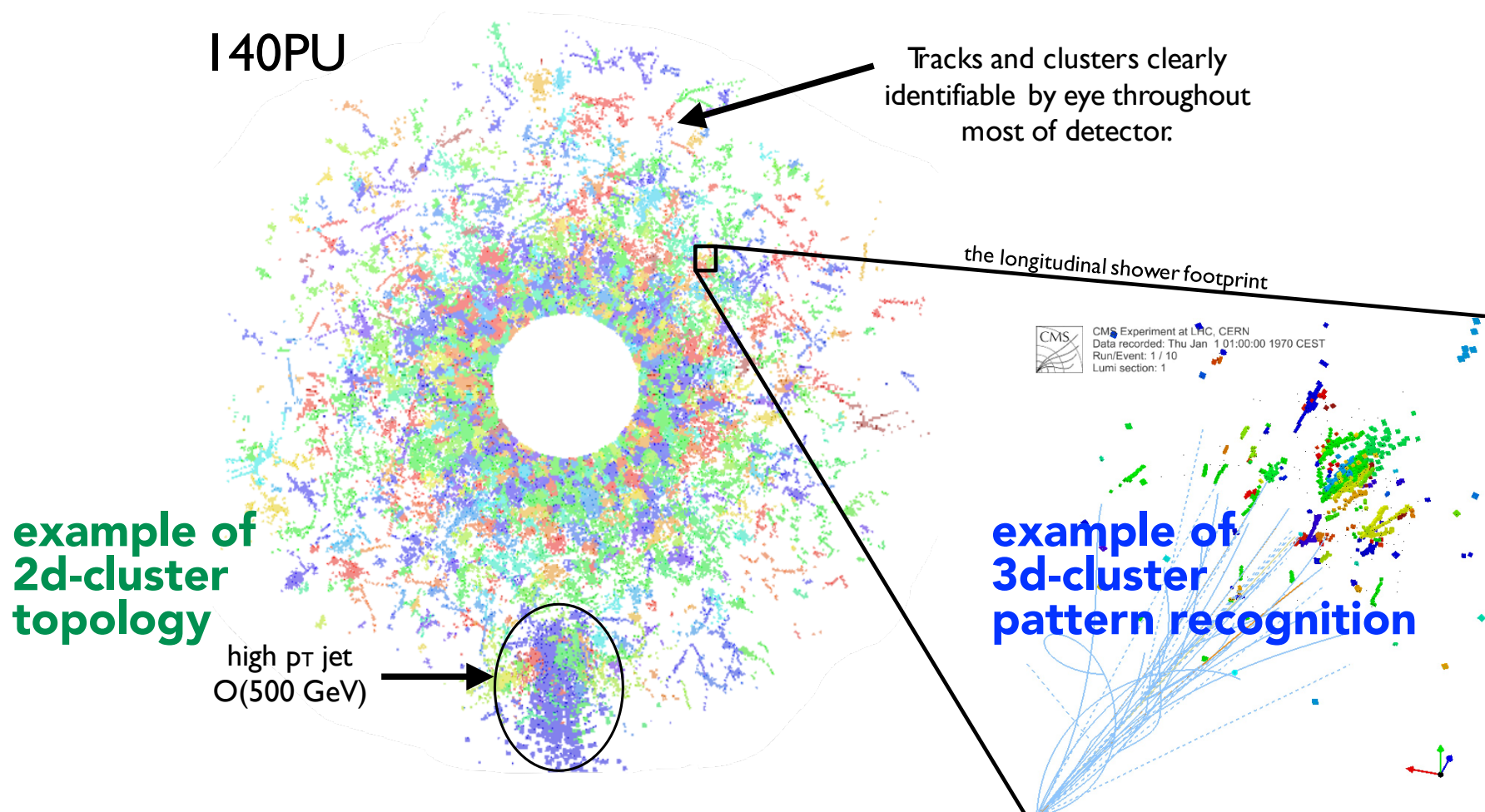**Building block of NOvA**

To 1 APD pixel

$\gamma$

typical charged particle path

L

6 cm   4 cm

Plane of vertical cells

Plane of horizontal cells

4 cm   6 cm

**Far Detector, on surface
14 kTons
896 readout planes
344,064 pixels**

Near   Proto

140PU

Tracks and clusters clearly identifiable by eye throughout most of detector.

the longitudinal shower footprint

CMS Experiment at LHC, CERN
Data recorded: Thu Jan 1 01:00:00 1970 CEST
Run/Event: 1 / 10
Lumi section: 1

**example of 2d-cluster topology**

high pᴛ jet O(500 GeV)

**example of 3d-cluster pattern recognition**

CMS week 2016    CERN June 20th - 24th

Privacy issue: not the focus today but probably deserves a plenary talk in some other conferences…